

1. *Naive Bayes Gaussian discriminant analysis.* Consider the generative model for (supervised) classification by assuming $\Pr[y_i = c] = \pi_c$ and $\mathbf{x}_i | y_i \sim \mathcal{N}(\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$. Without any additional assumptions, this is the Gaussian discriminant analysis (aka quadratic discriminant analysis) that we covered in class. We've discussed one special case when all the $\boldsymbol{\Sigma}_c$ matrices are the same, which leads to the linear discriminant analysis (LDA) model.

Here we make a different assumption: all the $\boldsymbol{\Sigma}_c$ matrices are diagonal (but not necessarily the same); for multivariate normals it means all variables are independent (conditioned on observing labels).

- (a) Derive the maximum likelihood estimate for the model parameters $\pi_c, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c, c = 1, \dots, k$.
- (b) Given a new data point \mathbf{x}_0 , explain how to predict \hat{y}_0 . Simplify the expression as much as possible.

Ans-1 =

Part(a) =

Bayes theorem Explains that -

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Hence we can say that,

$$P(A|B) \propto P(B|A)P(A)$$

In this case $A \Rightarrow [y_i = c]$
and $B = \mathbf{x}_i$

Then we can write that -

$$P(y_i=c|x_i) \propto P(x_i|y_i=c) P(y_i=c)$$

$\left\{ \begin{array}{l} \text{Here } P(y_i=c) = \pi \\ \text{and } P(x_i|y_i=c) = \mathcal{N}(\mu_c, \Sigma_c) \end{array} \right.$

Putting these values we get -

$$P(y_i=c|x_i) = \pi_c \mathcal{N}(\mu_c, \Sigma_c)$$

$$\Rightarrow \pi_c \det(2\pi\Sigma_c)^{-1/2} \exp\left(-\frac{1}{2}(x_i - \mu_c)^\top \Sigma_c^{-1} (x_i - \mu_c)\right)$$

maximum log likelihood function becomes

$$\Rightarrow \sum_{i=1}^n \log \pi_c - \frac{1}{2} \log 2\pi \det \Sigma_c - \frac{1}{2} (x_i - \mu_c)^\top \Sigma_c^{-1} (x_i - \mu_c)$$

① = Taking derivative of MLE function
with respect to π_c

$$\text{MLE} = \sum_{i=1}^n \log \pi_c \quad (\text{ignoring all constants})$$

construct lagrangian with $\lambda =$

$$L(\pi_c) = \sum_{i=1}^n \log \pi_c + \lambda \pi_c$$

$$\Rightarrow \frac{\partial L(\pi_c)}{\partial \pi_c} = \sum_{i=1}^n \frac{1}{\pi_c} + \lambda \cdot 1$$

(equate to 0)

$$\Rightarrow \sum_{i=1}^n \frac{1}{\pi_c} + \lambda = 0$$

$$\Rightarrow \pi_c = \frac{-n_c}{\lambda}$$

taking summation-

$$\sum_{c=1}^C \pi_c = -\frac{\sum_{c=1}^C (n_c)}{\lambda}$$

$$\boxed{\lambda = -\frac{N_c}{N}}$$

$$\boxed{\pi_c = \frac{\text{Total no. of Samples in class } c}{\text{Total no. of Samples in all classes}}}$$

② = calculating $\Sigma_C \Rightarrow$

partial derivative of MLE function

w.r.t. Σ_C we obtain =

$$\sum_{i=1}^n \left(\frac{\partial}{\partial \Sigma_C} \left(-\frac{1}{2} \log 2\pi \det \Sigma_C \right) \right)$$

$$\rightarrow \frac{\partial}{\partial \Sigma_C} \left[(x_i - \mu_C)^T \Sigma_C^{-1} (x_i - \mu_C) \right]$$

solving it we get =

$$\sum_{i=1}^n \left[\frac{1}{2} \Sigma_C^{-1} + \frac{1}{2} \Sigma_C^{-1} (x_i - \mu_C) (x_i - \mu_C)^T \Sigma_C^{-1} \right] = 0$$

$$\sum_{i=1}^n \left[\frac{1}{2} (I - \Sigma_C^{-1} (x_i - \mu_C) (x_i - \mu_C)^T) \Sigma_C^{-1} \right] = 0$$

multiplying both sides by $\Sigma_C \Rightarrow$

$$\sum_{i=1}^n \Sigma_C = \sum_{i=1}^n (x_i - \mu_C) (x_i - \mu_C)^T$$

Then Σ_C can be given as =

$$\Sigma_C = \frac{1}{N_C} \sum_{c=1}^C (x_i - \mu_c) (x_i - \mu_c)^T$$

③ = Obtaining μ_C -

Taking derivative of MLE w.r.t. μ_C .

$$\frac{\partial}{\partial \mu_C} \left[-\frac{1}{2} (x_i - \mu_C)^T \Sigma_C^{-1} (x_i - \mu_C) \right]$$

remaining all constant terms become = 0.

Σ_C^{-1} = diagonal \Rightarrow near symmetric

for Symmetric matrix =

$$\frac{\partial}{\partial A} A^T B A = 2BA$$

Hence above equation simplifies to -

$$\sum_{i=1}^n \Sigma_C^{-1} (x_i - \mu_c) = 0$$

which gives $\bar{w} =$

$$\boxed{\mu_c = \frac{1}{N_c} \sum_{i=1}^{N_c} x_i}$$

Part.b \Rightarrow Prediction (\hat{y}_0) for given value of x_0 .

$$\hat{y}_0 = f(x_0) = \operatorname{argmax}_c \pi_c \cdot P(y_i=c) \cdot P(x|y=c)$$

$$= \operatorname{argmax}_c \pi_c \cdot N(\mu_c, \Sigma_c)$$

$$= \operatorname{argmax}_c \left(\log \pi_c - \frac{1}{2} (x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c) \right)$$

remove all constant terms

$$= \operatorname{argmax}_c \left[-\frac{1}{2} x^T \Sigma_c^{-1} x + x^T \Sigma_c^{-1} \mu_c \right. \\ \left. - \frac{1}{2} \mu_c^T \Sigma_c^{-1} \mu_c + \log \pi_c \right]$$

2. Consider the lasso regularized k -class logistic regression problem

$$\underset{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_k}{\text{minimize}} \quad \frac{1}{n} \sum_{i=1}^n \left(\log \sum_{c=1}^k \exp(\boldsymbol{\theta}_c^\top \phi_i) - \boldsymbol{\theta}_{y_i}^\top \phi_i \right) + \frac{\lambda}{2} \sum_{c=1}^k \|\boldsymbol{\theta}_c\|_1.$$

Write the pseudocode of the proximal stochastic gradient descent algorithm for solving it.

Hint: Denote $f(\mathbf{z}) = \log \sum \exp(\mathbf{z})$, then $\nabla f(\mathbf{z}) = \frac{1}{\sum \exp(\mathbf{z})} \exp(\mathbf{z})$, where we overload the definition of $\exp(\cdot)$ for vector inputs by taking exponential element-wise. If $g(\mathbf{x}) = f(\mathbf{Ax})$, then $\nabla g(\mathbf{x}) = \mathbf{A}^\top \nabla f(\mathbf{Ax})$.

(a)= In order to minimize the objective function of the form \Rightarrow

$$\underset{\boldsymbol{\theta}}{\text{minimize}} \quad L(\boldsymbol{\theta}) + \lambda r(\boldsymbol{\theta})$$

$$\text{Here } L(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left(\log \sum_{c=1}^k \exp(\boldsymbol{\theta}_c^\top \phi_i) - \boldsymbol{\theta}_{y_i}^\top \phi_i \right)$$

$L(\boldsymbol{\theta})$ = Loss function

and $r(\boldsymbol{\theta})$ = Regularizer

$$r(\boldsymbol{\theta}) = \sum_{c=1}^k \|\boldsymbol{\theta}_c\|_1$$

In order to minimize this function and obtain the optimized values of Parameters $\boldsymbol{\theta}$ we have to use Proximal gradient descent.

Proximal gradient descent is given as =

$$\theta^{(t+1)} = \text{Prox}_{\gamma \lambda \rho}(\theta^+ - \gamma \nabla L(\theta^+))$$
$$= \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \|\theta - \theta^+ + \gamma \nabla L(\theta^+)\|^2 + \lambda \gamma \rho(\theta)$$

In order to obtain Proximal gradient we have to obtain gradient of loss function first i.e. $L(\theta)$

Pseudocode for calculating Proximal gradient descent =

① = initially assume parameter (θ) values.

② = obtain gradient of loss function $L(\theta)$.

Let's derive the gradient of $L(\theta)$.

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \left(\log \sum_{c=1}^k \exp(\theta_c^\top \phi_i) - \theta_j^\top \phi_i \right)$$

We can observe that -

$$f(z) = \log \sum \exp(z), \quad z = \theta_c^T \phi_i$$

Then we can write loss function as -

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n (f(z) - \theta_y^T \phi_i)$$

Taking gradient =

$$\nabla L(\theta) = \frac{1}{n} \sum_{i=1}^n (\nabla f(z) - (\nabla \theta_y^T) \cdot \phi_i - (\nabla \phi_i) \theta_y^T)$$

We are given that

$$\nabla f(z) = \frac{1}{\sum \exp(z)} \exp(z)$$

Then we can write \Rightarrow

$$\nabla L(\theta) = \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{\sum \exp(z)} - \nabla(\theta_y^T) \phi_i - (\nabla \phi_i) \theta_y^T \right]$$

Putting value of $z = \theta_c^T \phi_i$

we can write the loss function as =

$$\nabla L(\theta) = \frac{1}{n} \sum_{i=1}^n \left[\frac{\exp(\theta^T \phi_i)}{\sum_{j=1}^k \exp(\theta^T \phi_j)} - (\gamma \phi_i)^T \phi_i \right]$$

③ = Put the assumed value of parameters θ and value of one random sample from the data set (ϕ) in this loss function.

④ = Calculate value of $\theta^{(t+1)}$ for current iteration using proximal mapping formula.

$$\theta^{(t+1)} = \text{Prox}_{\gamma \lambda \sigma}(\theta^t - \gamma \nabla L(\theta^t))$$
$$= \underset{\theta}{\operatorname{arg\,min}} \frac{1}{2} \|\theta - \theta^t + \gamma \nabla L(\theta^t)\|^2 + \lambda \gamma \sigma(\theta)$$

⑤ = These steps will be iterated until the value of θ parameters converge.

3. Consider the latent variable model with the following probability distribution:

$$\Pr(y_i = c) = \pi_c, \quad \Pr(\mathbf{x}_i | y_i = c) \sim \text{Multi}(\mathbf{p}_c, L_i),$$

meaning that y_i is categorical, with k possible outcomes, and $\Pr(y_i = c) = \pi_c$; $(\mathbf{x}_i | y_i = c)$ follows a multinomial distribution by drawing from \mathbf{p}_c L_i times, i.e.,

$$p(\mathbf{x}_i | y_i = c) = \frac{L_i!}{\prod_{j=1}^d x_{ij}!} \prod_{j=1}^d p_{cj}^{x_{ij}}.$$

Given data samples $\mathbf{x}_1, \dots, \mathbf{x}_n$:

- (a) Write out the maximum likelihood formulation for estimating $\mathbf{p}_1, \dots, \mathbf{p}_k$, and $\boldsymbol{\pi}$. Simplify the objective function as much as possible.
- (b) Derive an expectation-maximization algorithm for approximately solving the aforementioned problem.
- (c) Implement this algorithm and try it on the 20 News Group data set with $k = 20$ (on the raw word-count data, without tf-idf preprocessing). Show the top 10 words in each cluster.

Ans-3 =

Part(a) =

given that $\Pr(y_i = c) = \pi_c$

and $\Pr(\mathbf{x}_i | y_i = c) = \text{Multi}(\mathbf{p}_c, L_i)$

also given that

$$\Pr(\mathbf{x}_i | y_i = c) = \frac{L_i!}{\prod_{j=1}^d x_{ij}!} \prod_{j=1}^d p_{cj}^{x_{ij}}$$

let Θ be $(\theta_1, \theta_2, \dots, \theta_k)$ be our shorthand unknown parameters.

max. log likelihood =

$$\begin{aligned} l_c(\theta; y) &= \sum_{i=1}^n \log p(x_i, y_i; \theta) \\ &= \sum_{i=1}^n \sum_{j=1}^k \pi(y_i=j) \log (\pi_j \phi(x_i; \theta_j)) \end{aligned}$$

Part(b) = algorithm =

- ① = initialize $\theta^{(0)}$ arbitrarily.
- ② = alternate b/w E & M step
until convergence.

$$a) = E\text{-Step} = q^{(t+1)} = \underset{q}{\operatorname{argmax}} L(q, \theta^t)$$

$$b) = M\text{-step} = \theta^{(t+1)} = \underset{\theta}{\operatorname{argmax}} L(q^{t+1}, \theta)$$

where $L(q, \theta)$ = Lower boundary term
auxiliary function

The solution is,

$$q_v(t++)(y) = p(y|x; \theta^+)$$

4. *Multidimensional scaling (MDS)*. MDS is another classical approach for unsupervised embedding, and to some extent relates to PCA.

The main idea of MDS is to embed each \mathbf{x}_i to \mathbf{y}_i so that the pair-wise distances are preserved as much as possible. This can be formulated as the following optimization problem

$$\underset{\mathbf{y}_1, \dots, \mathbf{y}_n}{\text{minimize}} \quad \sum_{i=1}^n \sum_{j=1}^{i-1} (\|\mathbf{x}_i - \mathbf{x}_j\| - \|\mathbf{y}_i - \mathbf{y}_j\|)^2.$$

There is no close-form solution for this formulation. A modified formulation called *classical scaling* is proposed:

$$\underset{\mathbf{y}_1, \dots, \mathbf{y}_n}{\text{minimize}} \quad \sum_{i=1}^n \sum_{j=1}^n (\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_j - \mathbf{y}_i^\top \mathbf{y}_j)^2, \quad (1)$$

where $\tilde{\mathbf{x}}_i = \mathbf{x}_i - (1/n) \sum_{j=1}^n \mathbf{x}_j$ is the centered data.

- (a) Use the Eckart-Young theorem to show that an optimal solution of (1) is to take the eigen-decomposition of the matrix $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} = \mathbf{U} \Lambda \mathbf{U}^\top$, where $\tilde{\mathbf{X}} = [\tilde{\mathbf{x}}_1 \cdots \tilde{\mathbf{x}}_n]$, keep the k largest eigenvalues in Λ and the corresponding columns in \mathbf{U} , and let $\mathbf{Y} = \Lambda^{1/2} \mathbf{U}^\top$; then each \mathbf{y}_i is the i th column of \mathbf{Y} .
- (b) Oftentimes one is directly given the pair-wise distance matrix \mathbf{D} , where

$$D_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|^2,$$

without explicitly given the data points $\mathbf{x}_1, \dots, \mathbf{x}_n$. Show that we can define the matrix

$$\mathbf{B} = \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right) \mathbf{D} \left(\mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top \right),$$

and replace $\tilde{\mathbf{x}}_i^\top \tilde{\mathbf{x}}_j$ with B_{ij} in formulation (1).

Part (a) =

formulated Problem for optimization =

$$\underset{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n}{\text{minimize}} \quad \sum_{i=1}^n \sum_{j=1}^{i-1} (\|\mathbf{x}_i - \mathbf{x}_j\| - \|\mathbf{y}_i - \mathbf{y}_j\|)^2$$

Classical Scaling - modified solution is -

$$\underset{y_1, y_2, \dots, y_n}{\text{minimize}} \sum_{i=1}^n \sum_{j=1}^n (\tilde{x}^T \tilde{x}_j - y_i^T y_j)^2$$

writing this equation in Matrix form -

$$\underset{y}{\text{minimize}} \| \tilde{x}^T \tilde{x} - y^T y \|$$

Eckart-Young theorem gives a low rank approximation, it gives solution to solve minimization problem. The cost function evaluates the fit between the data matrix and an approximating matrix which has reduced rank.

Mathematically Eckart-Young theorem can be written as -

for given matrices P and $Q \Rightarrow$

$$\| P - Q \| \geq \| P - P_Q \|$$

where (rank of matrix $Q \leq r$)
also P_r represents the best approximation
of P for rank r , considering the
largest k -eigen values and correspondi-
ng columns in SVD.

using the Eckart-Young theorem for
the above problem we can write =

$$\| \tilde{x}^T \tilde{x} - y^T y \| \geq \| \tilde{x}^T \tilde{x} - (\tilde{x}^T \tilde{x})_k \|$$

where we can say =

$(\tilde{x}^T \tilde{x})_k$ = approximation matrix
with rank k .

and rank of matrix $(y^T y) \leq k$.

Here we want to embed input vector
into a reduced dimension k .

we can say at optimal point =

$$Y^T Y = (\tilde{X}^T \tilde{X})_k \rightarrow \text{eqtn ①}$$

we can write SVD for $\tilde{X}^T \tilde{X}$ as =

$$\tilde{X}^T \tilde{X} = U \Lambda U^T$$

$\tilde{X}^T \tilde{X}$ is symmetric and Hence taking the largest k eigen values and corresponding columns in U we can write above equation as =

$$(\tilde{X}^T \tilde{X})_k = U_k \Lambda_k U_k^T$$

Putting the value of $(\tilde{X}^T \tilde{X})_k$ from equation ① for point of optimality =

we can write =

$$Y^T Y = U_k \Lambda_k U_k^T$$

Here $\mathbf{Y}^T \mathbf{Y}$ - represents a Positive Definite value.

if $\mathbf{Y}^T \mathbf{Y} \in (\text{tve, Definite value})$

Then

$$\boxed{\mathbf{Y} = \Lambda_k^{1/2} \mathbf{U}_k^T}$$

where $(\Lambda_k^{1/2})$ is also a positive non negative value

We have seen that $(\mathbf{Y} = \Lambda_k^{1/2} \mathbf{U}_k^T)$

after eigen value decomposition of $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$. we have got only k largest eigen values in Λ and corresponding columns in \mathbf{U} .

Hence \mathbf{Y} could be evaluated using relationship $\boxed{\mathbf{Y} = \Lambda_k^{1/2} \mathbf{U}_k^T}$

Proved

Part (b) =

The general conceptual process behind the idea of MDS is that →

- ① = First it transforms the squared distance matrix D to an inner product form.
- ② = Secondly it computes the eigen decomposition for this inner product form.

We are given the distance matrix D . for which

$$D_{ij} = \|x_i - x_j\|^2$$

$$D_{ij} = (x_i - x_j)^T (x_i - x_j)$$

$$D_{ij} = x_i^T x_i + x_j^T x_j - 2x_i^T x_j$$

→ equation ①

Let's apply the above mentioned process on given distance matrix $D =$

Let's assume that K is the inner product matrix of X .

$$K = X^T X \xrightarrow{\text{equation ②}}$$

with $K = \text{diag}(K_{ii}) \in \mathbb{R}^n$, Hence we can write equation ① by putting values of $X^T X$ from equation ② as

$$D = (d_{ij}^2) = K \cdot 1^T + 1 \cdot K^T - 2K$$

also we know that =

$$\tilde{x}_i = x_i - \text{mean}(X) = x_i - \frac{1}{n} X \cdot 1$$

$$\Rightarrow \tilde{X} = X - \frac{1}{n} X \cdot 1 \cdot 1^T \xrightarrow{\text{eqn ③}}$$

Then

$$\tilde{K} = \tilde{X}^T \tilde{X}$$

Putting value of \tilde{x} from equation ③
we get =

$$\tilde{K} = \tilde{x}^T \tilde{x}$$

$$= \left(x - \frac{1}{n} x \cdot 1 \cdot 1^T \right)^T \left(x - \frac{1}{n} x \cdot 1 \cdot 1^T \right)$$

$$= K - \frac{1}{n} K \cdot 1 \cdot 1^T - \frac{1}{n} 1 \cdot 1^T \cdot K + \frac{1}{n^2} \cdot 1 \cdot 1^T \cdot K \cdot 1 \cdot 1^T$$

let's assume that =

$$B = -\frac{1}{2} H \cdot D \cdot H^T$$

$$\text{where } H = I - \frac{1}{n} 1 \cdot 1^T$$

Here H is called the centering matrix.

$$\text{Then } B = -\frac{1}{2} H (K \cdot 1 \cdot 1^T + 1 \cdot K^T - 2K) \cdot H^T$$

Since,

$$\begin{aligned} K \cdot 1^T \cdot H^T &= K \cdot 1 \left(I - \frac{1}{n} 1 \cdot 1^T \right) \\ &= K \cdot 1 - K \left(\frac{1^T \cdot 1}{n} \right) \cdot 1 \\ &= 0 \end{aligned}$$

Hence, $H \cdot K 1 \cdot H^T = 0$, and

$$H \cdot 2K^T \cdot H^T = 0$$

Therefore,

$$\begin{aligned} B &= H \cdot K \cdot H^T \\ &= \left(I - \frac{1}{n} 1 \cdot 1^T \right) \cdot K \cdot \left(I - \frac{1}{n} 1 \cdot 1^T \right) \\ &= K - \frac{1}{n} 1 1^T K + \frac{1}{n} K \cdot 1 \cdot 1^T + \frac{1}{n^2} 1 \cdot 1^T (1^T \cdot K \cdot 1) \cdot 1 \\ &= K \end{aligned}$$

Hence $B = \tilde{K}$

and $\tilde{K} = \tilde{X}^T \tilde{X}$

we can write -

$$B = \tilde{X}^T \tilde{X}$$

Hence we have proved that
given a distance matrix $D = (d_{ij}^2)$
we can convert it to an inner
product matrix -

$$B = -\frac{1}{2} H D H^T$$

$$\text{where } H = \left(I - \frac{1}{n} \mathbf{1} \cdot \mathbf{1}^T \right)$$

is called centering matrix.

Problem formulation by replacing

$$B_{ij} = \tilde{x}_i^T \tilde{x}_j$$

Hence,

$$\text{minimize}_{y_1, y_2, \dots, y_n} \sum_{i=1}^n \sum_{j=1}^n (B_{ij} - y_i^T y_j)^2$$